

---

# Learning with a Wasserstein Loss

---

Charlie Frogner\*, Chiyuan Zhang\*, and Tomaso Poggio

Center for Brains, Minds and Machines  
McGovern Institute for Brain Research  
Massachusetts Institute of Technology  
frogner@mit.edu, chiyan@mit.edu, tp@ai.mit.edu

Hossein Mobahi  
MIT CSAIL

hmobahi@csail.mit.edu

Mauricio Araya-Polo

Shell International E & P, Inc.  
Mauricio.Araya@shell.com

## Abstract

Learning to predict multi-label outputs is challenging, but in many problems there is a natural metric on the outputs that can be used to improve predictions. In this paper we develop a loss function for multi-label learning, based on the Wasserstein distance. The Wasserstein distance provides a natural notion of dissimilarity for probability measures. Although optimizing with respect to the exact Wasserstein distance is costly, recent work has described a regularized approximation that is efficiently computed. We describe efficient learning algorithms based on this regularization, extending the Wasserstein loss from probability measures to unnormalized measures. We also describe a statistical learning bound for the loss and show connections with the total variation norm and the Jaccard index. The Wasserstein loss can encourage smoothness of the predictions with respect to a chosen metric on the output space. We demonstrate this property on a real-data tag prediction problem, using the Yahoo Flickr Creative Commons dataset, achieving superior performance over a baseline that doesn't use the metric.

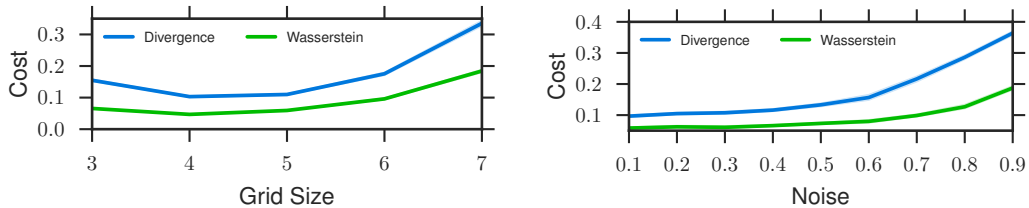
## 1 Introduction

We consider the problem of learning to predict a measure over a finite set. This problem includes many widely-used machine learning scenarios. For example, in multiclass classification, the set consists of the classes and a predicted distribution over classes is used to determine the top- $K$  most likely classes (as in the ImageNet Large Scale Visual Recognition Challenge [ILSVRC]) or to do subsequent inference (as with acoustic modeling in speech recognition). Another example is semantic segmentation [1], where the set consists of the pixel locations, and a segment can be modeled as a uniform measure supported on a subset.

In practice, many learning problems have natural similarity or metric structure on the output space. For example, in semantic segmentation, spatial adjacency between pixel locations provides a strong cue for similarity of their labels, due to contiguity of segmented regions. Such spatial adjacency can be captured, for example, by the Euclidean distance between the pixel locations. And in the ILSVRC image classification task, the output set comprises 1000 visual categories that are organized in a hierarchy, from which various semantic similarity measures are derived. Hierarchical structure in the label space is also prevalent in document categorization problems. In the following, we call the similarity structure in the label space the *ground metric* or *semantic similarity*.

---

\*Authors contributed equally.



(a) Cost vs. number of classes, averaged over different noise levels. (b) Cost vs. noise level (probability of mislabeling), averaged over different numbers of classes.

Figure 2: Confusion of near-equivalent classes degrades learning performance for a standard divergence-based loss. Incorporating semantic distance into the loss improves performance.

The presence of a ground metric can be taken into account when measuring the prediction performance. For example, confusing dogs with cats might be a more severe error than confusing breeds of dogs. Intuitively, a loss incorporating this metric should encourage the algorithm to favor predictions that are, if not completely accurate, at least semantically similar to the ground truth.

In this paper, we develop a loss function for multi-label learning that incorporates a metric on the output space by measuring the *Wasserstein distance* between a prediction and the target label, with respect to that metric. The Wasserstein distance is defined as the cost of the optimal transport plan for moving the mass in the predicted measure to match that in the target, and has been applied to a wide range of problems, including barycenter estimation [2], label propagation [3], and clustering [4]. To our knowledge, this paper represents the first use of the Wasserstein distance as a loss for supervised learning.



Figure 1: Semantically near-equivalent classes in ILSVRC

Incorporating an output metric into the loss can meaningfully impact learning performance. Take, for example, a multiclass classification problem containing semantically near-equivalent categories. Figure 1 shows such a case from the ILSVRC, in which the categories *Siberian husky* and *Eskimo dog* are nearly indistinguishable. Such categories can introduce noise in human-labeled data, as the labelers may fail to make fine distinctions between the categories. We simulate this problem by identifying the classes with points on a grid in the two-dimensional plane and randomly switching the labels to neighboring classes. We compare the standard multiclass logistic loss to the Wasserstein loss, and measure the prediction performance with the Euclidean distance between the predicted class and the true class. As shown in Figure 2, The prediction performance of both losses degrades as more labels are perturbed. Importantly, by incorporating the ground metric, the Wasserstein loss yields predictions that are closer to the ground truth, across all noise levels. Section D.1 of the Appendix describes the experiment in more detail.

The main contributions of this paper are as follows. We formulate the problem of learning with knowledge of the ground metric, and propose the Wasserstein loss as an alternative to traditional information divergence-based loss functions. Specifically, we focus on empirical risk minimization (ERM) with the Wasserstein loss, and describe efficient learning algorithms based on entropic regularization of the optimal transport problem. Moreover, we justify ERM with the Wasserstein loss by showing a statistical learning bound and we draw connections with existing measures of performance. Finally, we evaluate the proposed loss on both synthetic examples and a real-world image annotation problem, demonstrating benefits for incorporating an output metric into the loss.

## 2 Related work

Decomposable loss functions like KL Divergence and  $\ell_p$  distances are very popular for probabilistic [1] or vector-valued [5] predictions, as each component can be evaluated independently, often leading to simple and efficient algorithms. The idea of exploiting smoothness in the label space according to a prior metric has been explored in many different forms, including regularization [6] and post-processing with graphical models [7]. Optimal transport provides a natural distance for probability distributions over metric spaces. In [2, 8], the optimal transport is used to formulate the *Wasserstein Barycenter* as a probability distribution with minimum Wasserstein distance to a set of

given points on the probability simplex. [9] propagates histogram values on a graph by minimizing a Dirichlet energy induced by the optimal transport. The Wasserstein distance is also used to formulate a metric for comparing clusters in [10], as well as applied for image retrieval [11], contour matching [12], and many other problems that can be formulated as histogram matching [13]. However, to our knowledge, this is the first time it is used as a loss function in a discriminative learning framework. The closest work to this paper is a theoretical study [14] of an estimator that minimizes the optimal transport cost between the empirical distribution and the estimated distribution in the setting of statistical parameter estimation.

### 3 Learning with a Wasserstein loss

#### 3.1 Problem setup and notation

Consider the problem of learning a map from  $\mathcal{X} \in \mathbb{R}^{D_{\mathcal{X}}}$  to the space  $\mathcal{Y} = \mathbb{R}_+^K$  of measures over a finite set  $\mathcal{K}$  of size  $|\mathcal{K}| = K$ . Assume  $\mathcal{K}$  is a subset of a metric space with metric  $d_{\mathcal{K}}(\cdot, \cdot)$ .  $d_{\mathcal{K}}$  is called the *ground metric*, and it measures the semantic similarity in the label space. We perform learning over a hypothesis space  $\mathcal{H}$  of predictors  $h_{\theta} : \mathcal{X} \rightarrow \mathcal{Y}$ , parameterized by  $\theta \in \Theta$ .

In the standard statistical learning setting, we get an i.i.d. sequence of training examples  $S = ((x_1, y_1), \dots, (x_N, y_N))$ , sampled from an unknown joint distribution  $\mathcal{P}_{\mathcal{X} \times \mathcal{Y}}$ . Given a measure of performance (a.k.a. *risk*)  $\mathcal{E}(\cdot, \cdot)$ , the goal is to find the predictor  $h_{\theta} \in \mathcal{H}$  that minimizes the expected risk  $\mathbb{E}[\mathcal{E}(h_{\theta}(x), y)]$ . Typically  $\mathcal{E}(\cdot, \cdot)$  is difficult to optimize directly and the joint distribution  $\mathcal{P}_{\mathcal{X} \times \mathcal{Y}}$  is unknown, so learning is performed via *empirical risk minimization*. Specifically, we solve

$$\min_{h_{\theta} \in \mathcal{H}} \left\{ \hat{\mathbb{E}}_S[\ell(h_{\theta}(x), y)] = \frac{1}{N} \sum_{i=1}^N \ell(h_{\theta}(x_i), y_i) \right\}$$

with a loss function  $\ell(\cdot, \cdot)$  acting as a surrogate of  $\mathcal{E}(\cdot, \cdot)$ .

#### 3.2 Optimal transport and the exact Wasserstein loss

Information divergence-based loss functions are widely used in learning with probability-valued outputs. But along with other popular measures like Hellinger distance and  $\chi^2$  distance, these divergences are invariant to permutation of the elements in  $\mathcal{K}$ , ignoring any metric structure on  $\mathcal{K}$ .

Given a cost function  $c : \mathcal{K} \times \mathcal{K} \rightarrow \mathbb{R}$ , the *optimal transport* distance [15] measures the cheapest way to *transport* a probability measure  $\mu_1$  to match  $\mu_2$  with respect to  $c$ :

$$W_c(\mu_1, \mu_2) = \inf_{\gamma \in \Pi(\mu_1, \mu_2)} \int_{\mathcal{K} \times \mathcal{K}} c(\kappa_1, \kappa_2) \gamma(d\kappa_1, d\kappa_2) \quad (1)$$

where  $\Pi(\mu_1, \mu_2)$  is the set of joint probability measures on  $\mathcal{K} \times \mathcal{K}$  having  $\mu_1$  and  $\mu_2$  as marginals. An important case is when the cost is given by a metric  $d_{\mathcal{K}}(\cdot, \cdot)$  or its  $p$ -th power  $d_{\mathcal{K}}^p(\cdot, \cdot)$  with  $p \geq 1$ . In this case, they are called *Wasserstein distances* [16], also known as the *earth mover's distances* [11]. In this paper, we only work with discrete measures. In the case of probability measures, these are histograms in the simplex  $\Delta^{\mathcal{K}}$ .

When the ground truth  $y$  and the output of  $h$  both lie in the simplex  $\Delta^{\mathcal{K}}$ , we can define a Wasserstein loss at  $x$ .

**Definition 3.1** (Exact Wasserstein Loss). *For any  $h_{\theta} \in \mathcal{H}$ ,  $h_{\theta} : \mathcal{X} \rightarrow \Delta^{\mathcal{K}}$ , let  $h_{\theta}(\kappa|x) = h_{\theta}(x)_{\kappa}$  be the predicted value at element  $\kappa \in \mathcal{K}$ , given input  $x \in \mathcal{X}$ . Let  $y(\kappa)$  be the ground truth value for  $\kappa$  given by the corresponding label  $y$ . Then we define the Wasserstein loss as*

$$W_p^p(h(\cdot|x), y(\cdot)) = \inf_{T \in \Pi(h(x), y)} \langle T, M \rangle \quad (2)$$

where  $M \in \mathbb{R}_+^{K \times K}$  is the distance matrix  $M_{\kappa, \kappa'} = d_{\mathcal{K}}^p(\kappa, \kappa')$ , and the set of valid transport plans is

$$\Pi(h(x), y) = \{T \in \mathbb{R}_+^{K \times K} : T\mathbf{1} = h(x), T^{\top}\mathbf{1} = y\} \quad (3)$$

where  $\mathbf{1}$  is the all-one vector.

$W_p^p$  is the cost of the optimal plan for transporting the predicted mass distribution  $h(x)$  to match the target distribution  $y$ . The penalty increases as more mass is transported over longer distances, according to the ground metric  $M$ .

## 4 Efficient optimization

The Wasserstein loss (2) is a linear program and Lagrangian duality gives a means of computing descent direction with respect to  $h(x)$ . The dual LP of (2) is

$${}^dW_p^p(h(x), y) = \sup_{\alpha, \beta \in C_M} \alpha^\top h(x) + \beta^\top y, \quad C_M = \{(\alpha, \beta) \in \mathbb{R}^{K \times K} : \alpha_{\kappa} + \beta_{\kappa'} \leq M_{\kappa, \kappa'}\}. \quad (4)$$

As (2) is a linear program, at an optimum the values of the dual and the primal are equal (see, e.g. [17]), hence the dual optimal  $\alpha$  is a subgradient of the loss with respect to its first argument.

Computing  $\alpha$  is costly, as it entails solving a linear program with  $O(K^2)$  constraints, with  $K$  being the dimension of the output space. This cost can be prohibitive when optimizing by gradient descent.

### 4.1 Entropic regularization of optimal transport

Cuturi [18] proposes a smoothed transport objective that enables efficient approximation of both the transport matrix in (2) and the subgradient of the loss. [18] introduces an entropic regularization term that results in a strictly convex problem:

$${}^\lambda W_p^p(h(\cdot|x), y(\cdot)) = \inf_{T \in \Pi(h(x), y)} \langle T, M \rangle + \lambda H(T), \quad H(T) = - \sum_{\kappa, \kappa'} T_{\kappa, \kappa'} \log T_{\kappa, \kappa'}. \quad (5)$$

Importantly, the transport matrix that solves (5) is a *diagonal scaling* of a matrix  $\mathbf{K} = e^{-\lambda M - 1}$ :

$$T^* = \text{diag}(u) \mathbf{K} \text{diag}(v) \quad (6)$$

for  $u = e^{\lambda \alpha}$  and  $v = e^{\lambda \beta}$ , where  $\alpha$  and  $\beta$  are the Lagrangian dual variables for (5).

Identifying such a matrix subject to equality constraints on the row and column sums is exactly a *matrix balancing* problem, which is well-studied in numerical linear algebra and for which efficient iterative algorithms exist [19]. [18] and [2] use the well-known Sinkhorn-Knopp algorithm.

### 4.2 Extending smoothed transport to the learning setting

When the output vectors  $h(x)$  and  $y$  lie in the simplex, (5) can be used directly as a surrogate for (2). In this case,  $\alpha$  is a subgradient of the objective and can be obtained from the optimal scaling vector  $u$  as  $\alpha = \frac{1}{\lambda} \log u$ . Note that there is a translation ambiguity here: any upscaling of the vector  $u$  can be paired with a corresponding downscaling of the vector  $v$  without altering the matrix  $T^*$  (and vice versa). This means that  $\alpha$  is only defined up to a constant shift. In [2] the authors recommend choosing  $\alpha = \frac{1}{\lambda} \log u - \frac{1}{K\lambda} \log u^\top \mathbf{1}$  so that  $\alpha$  is tangent to the simplex.

For many learning problems, however, a normalized output assumption is unnatural. In image segmentation, for example, the target shape is not naturally represented as a histogram. And even when the prediction and the ground truth are constrained to the simplex, the observed label can be subject to noise that violates the constraint.

There is more than one way to generalize optimal transport to unnormalized measures. The objective we choose should deal effectively with the difference in total mass between  $h(x)$  and  $y$  while still being efficient to optimize.

### 4.3 Relaxed transport

We propose a novel relaxation that extends smoothed transport to unnormalized measures. By replacing the equality constraints on the transport marginals in (5) with soft penalties with respect to KL divergence, we get an unconstrained approximate transport problem. The resulting objective is:

$${}^{\lambda, \gamma_a, \gamma_b} W_{KL}(h(\cdot|x), y(\cdot)) = \min_{T \in \mathbb{R}_+^{K \times K}} \langle T, M \rangle + \lambda H(T) + \gamma_a \widetilde{\text{KL}}(T \mathbf{1} \| h(x)) + \gamma_b \widetilde{\text{KL}}(T^\top \mathbf{1} \| y) \quad (7)$$

where  $\widetilde{\text{KL}}(w \| z) = w^\top \log(w \oslash z) - \mathbf{1}^\top w + \mathbf{1}^\top z$  is the *generalized KL divergence* between  $w, z \in \mathbb{R}_+^K$ . Here  $\oslash$  represents element-wise division. As with the previous formulation, the optimal transport matrix with respect to (7) is a diagonal scaling of the matrix  $\mathbf{K}$ .

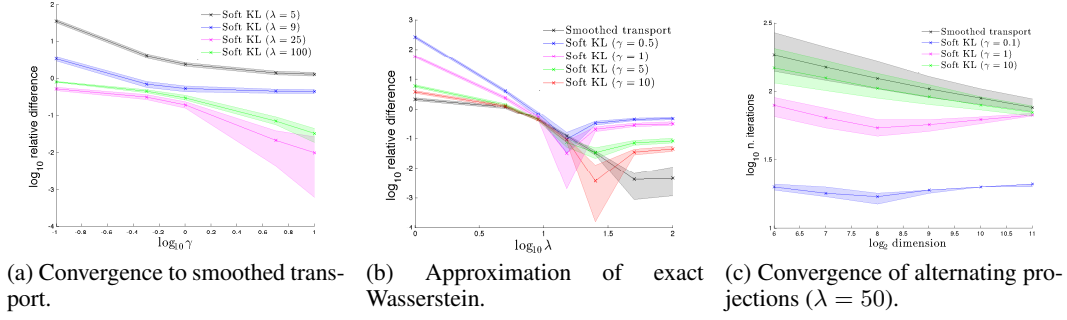


Figure 3: The relaxed transport problem (7) for unnormalized measures.

**Proposition 4.1.** *The transport matrix  $T^*$  optimizing (7) satisfies  $T^* = \text{diag}(u)\mathbf{K}\text{diag}(v)$ , where  $u = (h(x) \odot T^* \mathbf{1})^{\lambda\gamma_a}$ ,  $v = (y \odot (T^*)^\top \mathbf{1})^{\lambda\gamma_b}$ , and  $\mathbf{K} = e^{-\lambda M - 1}$ .*

And the optimal transport matrix is a fixed point for a Sinkhorn-like iteration.

**Proposition 4.2.**  *$T^* = \text{diag}(u)\mathbf{K}\text{diag}(v)$  optimizing (7) satisfies: i)  $u = h(x)^{\frac{\gamma_a\lambda}{\gamma_a\lambda+1}} \odot (\mathbf{K}v)^{-\frac{\gamma_a\lambda}{\gamma_a\lambda+1}}$ , and ii)  $v = y^{\frac{\gamma_b\lambda}{\gamma_b\lambda+1}} \odot (\mathbf{K}^\top u)^{-\frac{\gamma_b\lambda}{\gamma_b\lambda+1}}$ , where  $\odot$  represents element-wise multiplication.*

Unlike the previous formulation, (7) is unconstrained and differentiable with respect to  $h(x)$ . The gradient is given by  $\nabla_{h(x)} W_{KL}(h(\cdot|x), y(\cdot)) = \gamma_a (1 - T^* \mathbf{1} \odot h(x))$ .

When restricted to normalized measures, the relaxed problem (7) approximates smoothed transport (5). Figure 3a shows, for normalized  $h(x)$  and  $y$ , the relative distance between the values of (7) and (5)<sup>1</sup>. For  $\lambda$  large enough, (7) converges to (5) as  $\gamma_a$  and  $\gamma_b$  increase.

(7) also retains two properties of smoothed transport (5). Figure 3b shows that, for normalized outputs, the relaxed loss converges to the unregularized Wasserstein distance as  $\lambda$ ,  $\gamma_a$  and  $\gamma_b$  increase<sup>2</sup>. And Figure 3c shows that convergence of the iterations in (4.2) is nearly independent of the dimension  $K$  of the output space.

## 5 Properties of the Wasserstein loss

In this section, we study the statistical properties of learning with the exact Wasserstein loss (2) as well as connections with two standard measures. Full proofs can be found in the appendix.

### 5.1 Generalization error

Let  $S = ((x_1, y_1), \dots, (x_N, y_N))$  be i.i.d. samples and  $h_{\hat{\theta}}$  be the empirical risk minimizer

$$h_{\hat{\theta}} = \underset{h_{\theta} \in \mathcal{H}}{\text{argmin}} \left\{ \hat{\mathbb{E}}_S [W_p^p(h_{\theta}(\cdot|x), y)] = \frac{1}{N} \sum_{i=1}^N W_p^p(h_{\theta}(\cdot|x_i), y_i) \right\}.$$

Further assume  $\mathcal{H} = \mathfrak{s} \circ \mathcal{H}^o$  is the composition of a softmax  $\mathfrak{s}$  and a base hypothesis space  $\mathcal{H}^o$  of functions mapping into  $\mathbb{R}^K$ . The softmax layer outputs a prediction that lies in the simplex  $\Delta^K$ .

**Theorem 5.1.** *For  $p = 1$ , and any  $\delta > 0$ , with probability at least  $1 - \delta$ , it holds that*

$$\mathbb{E} [W_1^1(h_{\hat{\theta}}(\cdot|x), y)] \leq \inf_{h_{\theta} \in \mathcal{H}} \mathbb{E} [W_1^1(h_{\theta}(\cdot|x), y)] + 32KC_M \mathfrak{R}_N(\mathcal{H}^o) + 2C_M \sqrt{\frac{\log(1/\delta)}{2N}} \quad (8)$$

with the constant  $C_M = \max_{\kappa, \kappa'} M_{\kappa, \kappa'}$ .  $\mathfrak{R}_N(\mathcal{H}^o)$  is the Rademacher complexity [21] measuring the complexity of the hypothesis space  $\mathcal{H}^o$ .

<sup>1</sup>In figures 3a-c,  $h(x)$ ,  $y$  and  $M$  are generated as described in [18] section 5. In 3a-b,  $h(x)$  and  $y$  have dimension 256. In 3c, convergence is defined as in [18]. Shaded regions are 95% intervals.

<sup>2</sup>The unregularized Wasserstein distance was computed using `FastEMD` [20].

The Rademacher complexity  $\mathfrak{R}_N(\mathcal{H}^o)$  for commonly used models like neural networks and kernel machines [21] decays with the training set size. This theorem guarantees that the expected Wasserstein loss of the empirical risk minimizer approaches the best achievable loss for  $\mathcal{H}$ .

As an important special case, minimizing the empirical risk with Wasserstein loss is also good for multiclass classification. Let  $y = \mathbb{e}_\kappa$  be the ‘‘one-hot’’ encoded label vector for the groundtruth class.

**Proposition 5.2.** *In the multiclass classification setting, for  $p = 1$  and any  $\delta > 0$ , with probability at least  $1 - \delta$ , it holds that*

$$\mathbb{E}_{x, \kappa} [d_{\mathcal{K}}(\kappa_{\hat{\theta}}(x), \kappa)] \leq \inf_{h_{\theta} \in \mathcal{H}} K \mathbb{E}[W_1^1(h_{\theta}(x), y)] + 32K^2 C_M \mathfrak{R}_N(\mathcal{H}^o) + 2C_M K \sqrt{\frac{\log(1/\delta)}{2N}} \quad (9)$$

where the predictor is  $\kappa_{\hat{\theta}}(x) = \operatorname{argmax}_{\kappa} h_{\hat{\theta}}(\kappa|x)$ , with  $h_{\hat{\theta}}$  being the empirical risk minimizer.

Note that instead of the classification error  $\mathbb{E}_{x, \kappa}[\mathbb{1}\{\kappa_{\hat{\theta}}(x) \neq \kappa\}]$ , we actually get a bound on the expected semantic distance between the prediction and the groundtruth.

## 5.2 Connection with other standard measures

The special case in which no prior similarity is assumed between the points is captured by the 0-1 ground metric, defined by  $M_{\kappa, \kappa'}^{0-1} = \mathbb{1}_{\kappa \neq \kappa'} = 1 - \delta_{\kappa, \kappa'}$ . In this case, it is known that the Wasserstein distance reduces to the *total variation distance*  $\operatorname{TV}(\cdot, \cdot)$ :

**Proposition 5.3.** *For the 0-1 ground metric,  $\forall$  probability measures  $\mu, \nu$ ,  $W_{0-1}^1(\mu, \nu) = \operatorname{TV}(\mu, \nu)$ .*

The Wasserstein loss is also closely related to the *Jaccard index* [22], also known as intersection-over-union (IoU), which is a popular measure of performance in segmentation [23]. For two regions  $A$  and  $B$  in the image plane, the Jaccard index is defined as  $J(A, B) = |A \cap B|/|A \cup B|$ . The associated *Jaccard distance*  $d_J(A, B) = 1 - J(A, B)$  is a metric on the space of all finite sets [22]. If we treat each region  $A$  as a uniform probability distribution  $\mathbb{U}^A$  supported on  $A$ , then it holds that

**Proposition 5.4.** *The Wasserstein loss  $W_{0-1}^1$  is a proxy of  $d_J$  in the sense that for any  $0 \leq \varepsilon \leq 1$ ,  $W_{0-1}^1(\mathbb{U}^A, \mathbb{U}^B) \leq \varepsilon$  if  $d_J(A, B) \leq \varepsilon$ ; conversely,  $d_J(A, B) \leq 2\varepsilon$  if  $W_{0-1}^1(\mathbb{U}^A, \mathbb{U}^B) \leq \varepsilon$ .*

When the Euclidean distance in the image plane is used as the ground metric, the general Wasserstein loss  $W_p^p$  is still a surrogate of  $d_J$ :

**Corollary 5.5.** *For any  $0 \leq \varepsilon \leq 1$ , and  $p \geq 1$ , under the ground metric  $d(\kappa, \kappa') = \|\kappa - \kappa'\|_p^p$  over the set of pixel coordinates,  $W_p^p(\mathbb{U}^A, \mathbb{U}^B) \leq \varepsilon$  implies that  $d_J(A, B) \leq 2\varepsilon$ .*

Unlike  $W_{0-1}^1$ ,  $W_p^p$  is stronger than  $d_J$  because it ensures not only that the incorrectly predicted region is small, but also that it is not far away. The connection with the Jaccard distance can also be characterized for the case of non-uniform distributions. See section C.4 in the Appendix for details.

## 6 Empirical study

### 6.1 Impact of the ground metric

In this section, we show that the Wasserstein loss encourages smoothness with respect to an artificial metric on the MNIST handwritten digit dataset. This is a multi-class classification problem with output dimensions corresponding to the 10 digits, and we apply a ground metric  $d_p(\kappa, \kappa') = |\kappa - \kappa'|^p$ , where  $\kappa, \kappa' \in \{0, \dots, 9\}$  and  $p \in [0, \infty)$ . This metric encourages the recognized digit to be *numerically* close to the true one. We train a model independently for each value of  $p$  and plot the average predicted probabilities of the different digits on the test set in Figure 4.

Note that as  $p \rightarrow 0$ , the metric approaches the 0-1 metric  $d_0(\kappa, \kappa') = \mathbb{1}_{\kappa \neq \kappa'}$ , which treats all incorrect digits as being equally unfavorable. In this case, as can be seen in the figure, the predicted probability of the true digit goes to 1 while the probability for all other digits goes to 0. As  $p$  increases, the predictions become more evenly distributed over the neighboring digits, converging to a uniform distribution as  $p \rightarrow \infty$ <sup>3</sup>.

<sup>3</sup>To avoid numerical issues, we scale down the ground metric such that all of the distance values are in the interval  $[0, 1)$ .

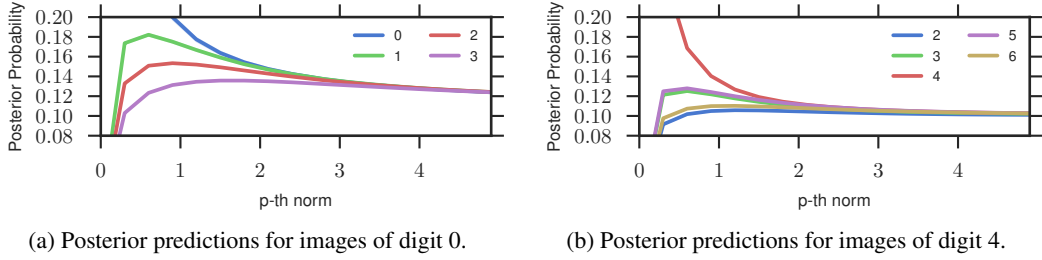


Figure 4: MNIST example. Each curve shows the predicted probability for one digit, for models trained with different  $p$  values for the ground metric.

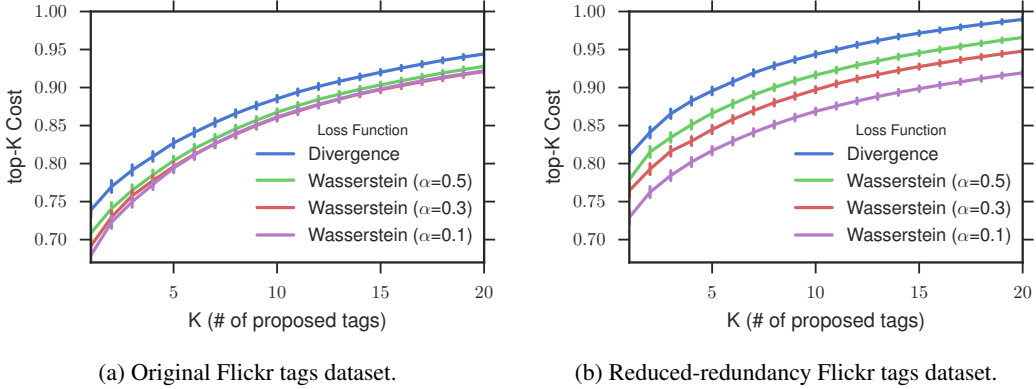


Figure 5: Top-K cost comparison of the proposed loss (Wasserstein) and the baseline (Divergence).

## 6.2 Flickr tag prediction

We apply the Wasserstein loss to a real world multi-label learning problem, using the recently released Yahoo/Flickr Creative Commons 100M dataset [24]. Our goal is *tag prediction*: we select 1000 descriptive tags along with two random sets of 10,000 images each, associated with these tags, for training and testing. We derive a distance metric between tags by using `word2vec` [25] to embed the tags as unit vectors, then taking their Euclidean distances. To extract image features we use `MatConvNet` [26]. Note that the set of tags is highly redundant and often many semantically equivalent or similar tags can apply to an image. The images are also incompletely tagged, as different users may prefer different tags. We therefore measure the prediction performance by the *top-K cost*, defined as  $C_K = 1/K \sum_{k=1}^K \min_j d_K(\hat{\kappa}_k, \kappa_j)$ , where  $\{\kappa_j\}$  is the set of groundtruth tags, and  $\{\hat{\kappa}_k\}$  are the tags with highest predicted probability.

We find that a linear combination of the Wasserstein loss  $W_p^p$  and a KL divergence-based loss yields the best prediction results. Specifically, we train a linear model by minimizing  $W_p^p + \alpha \text{KL}$  on the training set, where  $\alpha$  controls the relative weight of KL. Figure 5a shows the top-K cost on the test set for the combined loss and the baseline KL loss. We additionally create a second dataset by removing redundant labels from the original dataset: this simulates the potentially more difficult case in which a single user tags each image, by selecting one tag to apply from amongst each cluster of applicable, semantically similar tags. Figure 3b shows that performance for both algorithms decreases on the harder dataset, while the combined Wasserstein loss continues to outperform the baseline.

In Figure 6, we show the effect on performance of varying the weight  $\alpha$  on the KL loss. We observe that the optimum of the top-K cost is achieved when the Wasserstein loss is weighted more heavily than at the optimum of the AUC. This is consistent with a semantic smoothing effect of Wasserstein, which during training will favor mispredictions that are semantically similar to the ground truth, sometimes at the cost of lower AUC<sup>4</sup>. We finally show two selected images from the test set in

<sup>4</sup>The Wasserstein loss can achieve a similar trade-off alone as discussed in Section 6.1. However, the achievable range is usually limited by numerical stability when dealing with large values of the metric. In practice it is often easier to implement the trade-off by combining with a KL loss.



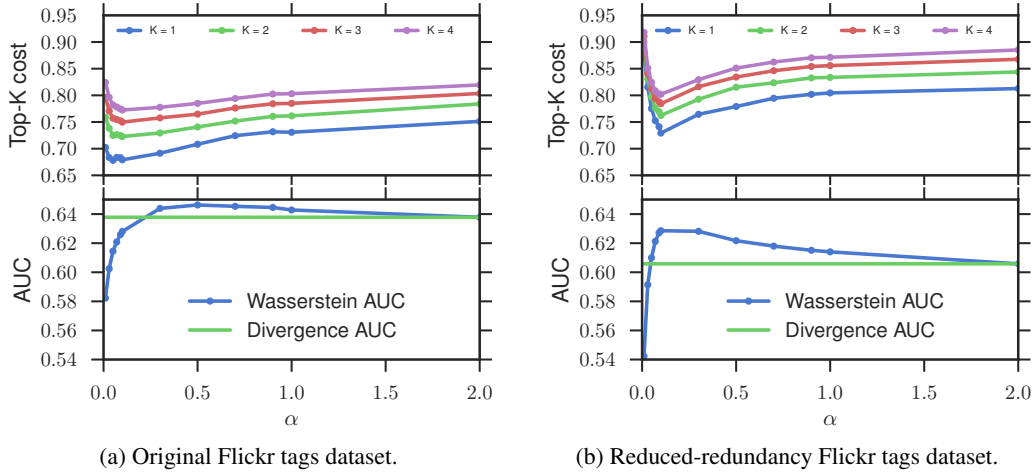


Figure 6: Trade-off between semantic smoothness and maximum likelihood.



(a) **Flickr user tags:** street, parade, dragon; **our proposals:** people, protest, parade; **baseline proposals:** music, car, band.



(b) **Flickr user tags:** water, boat, reflection, sunshine; **our proposals:** water, river, lake, summer; **baseline proposals:** river, water, club, nature.

Figure 7: Examples of images in the Flickr dataset. We show the groundtruth tags and as well as tags proposed by our algorithm and the baseline.

Figure 7. These illustrate cases in which both algorithms make predictions that are semantically relevant, despite overlapping very little with the ground truth. The image on the left shows semantically irrelevant errors made by the baseline algorithm. More examples can be found in the appendix.

## 7 Conclusions and future work

In this paper we have described a loss function for learning to predict a measure over a finite set, based on the Wasserstein distance. Optimizing with respect to the exact Wasserstein loss is computationally costly and we describe efficient algorithms based on entropic regularization, for learning both normalized and unnormalized measures. We have also described a statistical learning bound for the loss and shown connections with both the total variation norm and the Jaccard index. The Wasserstein loss can encourage smoothness of the predictions with respect to a chosen metric on the output space, and we demonstrate this property on a real-data tag prediction problem, achieving superior performance over a baseline that doesn't incorporate the metric.

An interesting direction for future work may be to explore the connection between the Wasserstein loss and Markov random fields, as the latter are often used to encourage smoothness of predictions, via inference at prediction time.



## References

- [1] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *CVPR (to appear)*, 2015.
- [2] Marco Cuturi and Arnaud Doucet. Fast Computation of Wasserstein Barycenters. *ICML*, 2014.
- [3] Justin Solomon, Raif M Rustamov, Leonidas J Guibas, and Adrian Butscher. Wasserstein Propagation for Semi-Supervised Learning. *ICML*, pages 306–314, 2014.
- [4] Michael H Coen, M Hidayath Ansari, and Nathanael Fillmore. Comparing Clusterings in Space. *ICML*, pages 231–238, 2010.
- [5] Lorenzo Rosasco Mauricio A. Alvarez and Neil D. Lawrence. Kernels for vector-valued functions: A review. *Foundations and Trends in Machine Learning*, 4(3):195–266, 2011.
- [6] Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, 1992.
- [7] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015.
- [8] Marco Cuturi, Gabriel Peyré, and Antoine Rolet. A Smoothed Dual Approach for Variational Wasserstein Problems. *arXiv.org*, March 2015.
- [9] Justin Solomon, Raif Rustamov, Leonidas Guibas, and Adrian Butscher. Wasserstein propagation for semi-supervised learning. In *ICML*, 2014.
- [10] Michael Coen, Hidayath Ansari, and Nathanael Fillmore. Comparing clusterings in space. In *ICML*, 2010.
- [11] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. The earth mover’s distance as a metric for image retrieval. *IJCV*, 40(2):99–121, 2000.
- [12] Kristen Grauman and Trevor Darrell. Fast contour matching using approximate earth mover’s distance. In *CVPR*, 2004.
- [13] S Shirdhonkar and D W Jacobs. Approximate earth mover’s distance in linear time. In *CVPR*, 2008.
- [14] Federico Bassetti, Antonella Bodini, and Eugenio Regazzini. On minimum kantorovich distance estimators. *Stat. Probab. Lett.*, 76(12):1298–1302, 1 July 2006.
- [15] Cédric Villani. *Optimal Transport: Old and New*. Springer Berlin Heidelberg, 2008.
- [16] Vladimir I Bogachev and Aleksandr V Kolesnikov. The Monge-Kantorovich problem: achievements, connections, and perspectives. *Russian Math. Surveys*, 67(5):785, 10 2012.
- [17] Dimitris Bertsimas, John N. Tsitsiklis, and John Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific, Boston, third printing edition, 1997.
- [18] Marco Cuturi. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. *NIPS*, 2013.
- [19] Philip A Knight and Daniel Ruiz. A fast algorithm for matrix balancing. *IMA Journal of Numerical Analysis*, 33(3):drs019–1047, October 2012.
- [20] Ofir Pele and Michael Werman. Fast and robust Earth Mover’s Distances. *ICCV*, pages 460–467, 2009.
- [21] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *JMLR*, 3:463–482, March 2003.
- [22] Michael Levandowsky and David Winter. Distance between sets. *Nature*, 234(5323):34–35, 1971.
- [23] Sebastian Nowozin. Optimal Decisions from Probabilistic Models: The Intersection-over-Union Case. *CVPR*, pages 548–555, 2014.
- [24] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. The new data and new challenges in multimedia research. *arXiv preprint arXiv:1503.01817*, 2015.
- [25] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [26] A. Vedaldi and K. Lenc. MatConvNet – Convolutional Neural Networks for MATLAB. *CoRR*, abs/1412.4564, 2014.
- [27] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Classics in Mathematics. Springer Berlin Heidelberg, 2011.
- [28] Clark R. Givens and Rae Michael Shortt. A class of wasserstein metrics for probability distributions. *Michigan Math. J.*, 31(2):231–240, 1984.

## A Relaxed transport

Equation (7) gives the relaxed transport objective as

$$^{\lambda, \gamma_a, \gamma_b} W_{KL}(h(\cdot|x), y(\cdot)) = \min_{T \in \mathbb{R}_+^{K \times K}} \langle T, M \rangle + \lambda H(T) + \gamma_a \widetilde{\mathbf{KL}}(T \mathbf{1} \| h(x)) + \gamma_b \widetilde{\mathbf{KL}}(T^\top \mathbf{1} \| y)$$

with  $\widetilde{\mathbf{KL}}(w \| z) = w^\top \log(w \odot z) - \mathbf{1}^\top w + \mathbf{1}^\top z$ .

*Proof of Proposition 4.1.* The first order condition for  $T^*$  optimizing (7) is

$$\begin{aligned} M_{ij} + \frac{1}{\lambda} (\log T_{ij}^* + 1) + \gamma_a (\log T^* \mathbf{1} \odot h(x))_i + \gamma_b (\log (T^*)^\top \mathbf{1} \odot y)_j &= 0. \\ \Rightarrow \log T_{ij}^* + \gamma_a \lambda \log (T^* \mathbf{1})_i + \gamma_b \lambda \log ((T^*)^\top \mathbf{1})_j &= -\lambda M_{ij} + \gamma_a \lambda \log h(x)_i + \gamma_b \lambda \log y_j - 1 \\ \Rightarrow T_{ij}^* (T^* \mathbf{1})^{\gamma_a \lambda} ((T^*)^\top \mathbf{1})^{\gamma_b \lambda} &= \exp(-\lambda M_{ij} + \gamma_a \lambda \log h(x)_i + \gamma_b \lambda \log y_j - 1) \\ \Rightarrow T_{ij}^* &= (h(x) \odot T^* \mathbf{1})_i^{\gamma_a \lambda} (y \odot (T^*)^\top \mathbf{1})_j^{\gamma_b \lambda} \exp(-\lambda M_{ij} - 1) \end{aligned}$$

Hence  $T^*$  (if it exists) is a diagonal scaling of  $\mathbf{K} = \exp(-\lambda M - 1)$ . □

*Proof of Proposition 4.2.* Let  $u = (h(x) \odot T^* \mathbf{1})^{\gamma_a \lambda}$  and  $v = (y \odot (T^*)^\top \mathbf{1})^{\gamma_b \lambda}$ , so  $T^* = \text{diag}(u) \mathbf{K} \text{diag}(v)$ . We have

$$\begin{aligned} T^* \mathbf{1} &= \text{diag}(u) \mathbf{K} v \\ \Rightarrow (T^* \mathbf{1})^{\gamma_a \lambda + 1} \odot h(x)^{\gamma_a \lambda} &= \mathbf{K} v \end{aligned}$$

where we substituted the expression for  $u$ . Re-writing  $T^* \mathbf{1}$ ,

$$\begin{aligned} (\text{diag}(u) \mathbf{K} v)^{\gamma_a \lambda + 1} &= \text{diag}(h(x)^{\gamma_a \lambda}) \mathbf{K} v \\ \Rightarrow u^{\gamma_a \lambda + 1} &= h(x)^{\gamma_a \lambda} \odot (\mathbf{K} v)^{-\gamma_a \lambda} \\ \Rightarrow u &= h(x)^{\frac{\gamma_a \lambda}{\gamma_a \lambda + 1}} \odot (\mathbf{K} v)^{-\frac{\gamma_a \lambda}{\gamma_a \lambda + 1}}. \end{aligned}$$

A symmetric argument shows that  $v = y^{\frac{\gamma_b \lambda}{\gamma_b \lambda + 1}} \odot (\mathbf{K}^\top u)^{-\frac{\gamma_b \lambda}{\gamma_b \lambda + 1}}$ . □

## B Statistical Learning Bounds

We establish the proof of Theorem 5.1 in this section. For simpler notation, for a sequence  $S = ((x_1, y_1), \dots, (x_N, y_N))$  of i.i.d. training samples, we denote the empirical risk  $\hat{R}_S$  and risk  $R$  as

$$\hat{R}_S(h_\theta) = \hat{\mathbb{E}}_S [W_p^p(h_\theta(\cdot|x), y(\cdot))], \quad R(h_\theta) = \mathbb{E} [W_p^p(h_\theta(\cdot|x), y(\cdot))] \quad (10)$$

**Lemma B.1.** Let  $h_{\hat{\theta}}, h_{\theta^*} \in \mathcal{H}$  be the minimizer of the empirical risk  $\hat{R}_S$  and expected risk  $R$ , respectively. Then

$$R(h_{\hat{\theta}}) \leq R(h_{\theta^*}) + 2 \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)|$$

*Proof.* By the optimality of  $h_{\hat{\theta}}$  for  $\hat{R}_S$ ,

$$\begin{aligned} R(h_{\hat{\theta}}) - R(h_{\theta^*}) &= R(h_{\hat{\theta}}) - \hat{R}_S(h_{\hat{\theta}}) + \hat{R}_S(h_{\hat{\theta}}) - R(h_{\theta^*}) \\ &\leq R(h_{\hat{\theta}}) - \hat{R}_S(h_{\hat{\theta}}) + \hat{R}_S(h_{\theta^*}) - R(h_{\theta^*}) \\ &\leq 2 \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \end{aligned}$$

□

Therefore, to bound the risk for  $h_{\hat{\theta}}$ , we need to establish uniform concentration bounds for the Wasserstein loss. Towards that goal, we define a space of loss functions induced by the hypothesis space  $\mathcal{H}$  as

$$\mathcal{L} = \{\ell_\theta : (x, y) \mapsto W_p^p(h_\theta(\cdot|x), y(\cdot)) : h_\theta \in \mathcal{H}\} \quad (11)$$

The uniform concentration will depends on the “complexity” of  $\mathcal{L}$ , which is measured by the empirical *Rademacher complexity* defined below.

**Definition B.2** (Rademacher Complexity [21]). *Let  $\mathcal{G}$  be a family of mapping from  $\mathcal{Z}$  to  $\mathbb{R}$ , and  $S = (z_1, \dots, z_N)$  a fixed sample from  $\mathcal{Z}$ . The empirical Rademacher complexity of  $\mathcal{G}$  with respect to  $S$  is defined as*

$$\hat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \frac{1}{N} \sum_{i=1}^n \sigma_i g(z_i) \right] \quad (12)$$

where  $\sigma = (\sigma_1, \dots, \sigma_N)$ , with  $\sigma_i$ 's independent uniform random variables taking values in  $\{+1, -1\}$ .  $\sigma_i$ 's are called the *Rademacher random variables*. The Rademacher complexity is defined by taking expectation with respect to the samples  $S$ ,

$$\mathfrak{R}_N(\mathcal{G}) = \mathbb{E}_S [\hat{\mathfrak{R}}_S(\mathcal{G})] \quad (13)$$

**Theorem B.3.** *For any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following holds for all  $\ell_\theta \in \mathcal{L}$ ,*

$$\mathbb{E}[\ell_\theta] - \hat{\mathbb{E}}_S[\ell_\theta] \leq 2\mathfrak{R}_N(\mathcal{L}) + \sqrt{\frac{C_M^2 \log(1/\delta)}{2N}} \quad (14)$$

with the constant  $C_M = \max_{\kappa, \kappa'} M_{\kappa, \kappa'}$ .

By the definition of  $\mathcal{L}$ ,  $\mathbb{E}[\ell_\theta] = R(h_\theta)$  and  $\hat{\mathbb{E}}_S[\ell_\theta] = \hat{R}_S[h_\theta]$ . Therefore, this theorem provides a uniform control for the deviation of the empirical risk from the risk.

**Theorem B.4** (McDiarmid's Inequality). *Let  $S = \{X_1, \dots, X_N\} \subset \mathcal{X}$  be  $N$  i.i.d. random variables. Assume there exists  $C > 0$  such that  $f : \mathcal{X}^N \rightarrow \mathbb{R}$  satisfies the following stability condition*

$$|f(x_1, \dots, x_i, \dots, x_N) - f(x_1, \dots, x'_i, \dots, x_N)| \leq C \quad (15)$$

for all  $i = 1, \dots, N$  and any  $x_1, \dots, x_N, x'_i \in \mathcal{X}$ . Then for any  $\varepsilon > 0$ , denoting  $f(X_1, \dots, X_N)$  by  $f(S)$ , it holds that

$$\mathbb{P}(f(S) - \mathbb{E}[f(S)] \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{NC^2}\right) \quad (16)$$

**Lemma B.5.** *Let the constant  $C_M = \max_{\kappa, \kappa'} M_{\kappa, \kappa'}$ , then  $0 \leq W_p^p(\cdot, \cdot) \leq C_M$ .*

*Proof.* For any  $h(\cdot|x)$  and  $y(\cdot)$ , let  $T^* \in \Pi(h(x), y)$  be the optimal transport plan that solves (2), then

$$W_p^p(h(x), y) = \langle T^*, M \rangle \leq C_M \sum_{\kappa, \kappa'} T_{\kappa, \kappa'} = C_M$$

□

*Proof of Theorem B.3.* For any  $\ell_\theta \in \mathcal{L}$ , note the empirical expectation is the empirical risk of the corresponding  $h_\theta$ :

$$\hat{\mathbb{E}}_S[\ell_\theta] = \frac{1}{N} \sum_{i=1}^N \ell_\theta(x_i, y_i) = \frac{1}{N} \sum_{i=1}^N W_p^p(h_\theta(\cdot|x_i), y_i(\cdot)) = \hat{R}_S(h_\theta)$$

Similarly,  $\mathbb{E}[\ell_\theta] = R(h_\theta)$ . Let

$$\Phi(S) = \sup_{\ell \in \mathcal{L}} \mathbb{E}[\ell] - \hat{\mathbb{E}}_S[\ell] \quad (17)$$

Let  $S'$  be  $S$  with the  $i$ -th sample replaced by  $(x'_i, y'_i)$ , by Lemma B.5, it holds that

$$\Phi(S) - \Phi(S') \leq \sup_{\ell \in \mathcal{L}} \hat{\mathbb{E}}_{S'}[\ell] - \hat{\mathbb{E}}_S[\ell] = \sup_{h_\theta \in \mathcal{H}} \frac{W_p^p(h_\theta(x'_i), y'_i) - W_p^p(h_\theta(x_i), y_i)}{N} \leq \frac{C_M}{N}$$

Similarly, we can show  $\Phi(S') - \Phi(S) \leq C_M/N$ , thus  $|\Phi(S') - \Phi(S)| \leq C_M/N$ . By Theorem B.4, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , it holds that

$$\Phi(S) \leq \mathbb{E}[\Phi(S)] + \sqrt{\frac{C_M^2 \log(1/\delta)}{2N}} \quad (18)$$

To bound  $\mathbb{E}[\Phi(S)]$ , by Jensen's inequality,

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_{\ell \in \mathcal{L}} \mathbb{E}[\ell] - \hat{\mathbb{E}}_S[\ell] \right] = \mathbb{E}_S \left[ \sup_{\ell \in \mathcal{L}} \mathbb{E}_{S'} \left[ \hat{\mathbb{E}}_{S'}[\ell] - \hat{\mathbb{E}}_S[\ell] \right] \right] \leq \mathbb{E}_{S, S'} \left[ \sup_{\ell \in \mathcal{L}} \hat{E}_{S'}[\ell] - \hat{E}_S[\ell] \right]$$

Here  $S'$  is another sequence of i.i.d. samples, usually called *ghost samples*, that is only used for analysis. Now we introduce the Rademacher variables  $\sigma_i$ , since the role of  $S$  and  $S'$  are completely symmetric, it follows

$$\begin{aligned} \mathbb{E}_S[\Phi(S)] &\leq \mathbb{E}_{S, S', \sigma} \left[ \sup_{\ell \in \mathcal{L}} \frac{1}{N} \sum_{i=1}^N \sigma_i (\ell(x'_i, y'_i) - \ell(x_i, y_i)) \right] \\ &\leq \mathbb{E}_{S', \sigma} \left[ \sup_{\ell \in \mathcal{L}} \frac{1}{N} \sum_{i=1}^N \sigma_i \ell(x'_i, y'_i) \right] + \mathbb{E}_{S, \sigma} \left[ \sup_{\ell \in \mathcal{L}} \frac{1}{N} \sum_{i=1}^N -\sigma_i \ell(x_i, y_i) \right] \\ &= \mathbb{E}_S \left[ \hat{\mathfrak{R}}_S(\mathcal{L}) \right] + \mathbb{E}_{S'} \left[ \hat{\mathfrak{R}}_{S'}(\mathcal{L}) \right] \\ &= 2\mathfrak{R}_N(\mathcal{L}) \end{aligned}$$

The conclusion follows by combining (17) and (18).  $\square$

To finish the proof of Theorem 5.1, we combine Lemma B.1 and Theorem B.3, and relate  $\mathfrak{R}_N(\mathcal{L})$  to  $\mathfrak{R}_N(\mathcal{H})$  via the following generalized Talagrand's lemma [27].

**Lemma B.6.** *Let  $\mathcal{F}$  be a class of real functions, and  $\mathcal{H} \subset \mathcal{F} = \mathcal{F}_1 \times \dots \times \mathcal{F}_K$  be a  $K$ -valued function class. If  $\mathfrak{m} : \mathbb{R}^K \rightarrow \mathbb{R}$  is a  $L_m$ -Lipschitz function and  $\mathfrak{m}(0) = 0$ , then  $\hat{\mathfrak{R}}_S(\mathfrak{m} \circ \mathcal{H}) \leq 2L_m \sum_{k=1}^K \hat{\mathfrak{R}}_S(\mathcal{F}_k)$ .*

**Theorem B.7** (Theorem 6.15 of [15]). *Let  $\mu$  and  $\nu$  be two probability measures on a Polish space  $(\mathcal{K}, d_{\mathcal{K}})$ . Let  $p \in [1, \infty)$  and  $\kappa_0 \in \mathcal{K}$ . Then*

$$W_p(\mu, \nu) \leq 2^{1/p'} \left( \int_{\mathcal{K}} d_{\mathcal{K}}(\kappa_0, \kappa) d|\mu - \nu|(\kappa) \right)^{1/p}, \quad \frac{1}{p} + \frac{1}{p'} = 1 \quad (19)$$

**Corollary B.8.** *The Wasserstein loss is Lipschitz continuous in the sense that for any  $h_{\theta} \in \mathcal{H}$ , and any  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ ,*

$$W_p^p(h_{\theta}(\cdot|x), y) \leq 2^{p-1} C_M \sum_{\kappa \in \mathcal{K}} |h_{\theta}(\kappa|x) - y(\kappa)| \quad (20)$$

*In particular, when  $p = 1$ , we have*

$$W_1^1(h_{\theta}(\cdot|x), y) \leq C_M \sum_{\kappa \in \mathcal{K}} |h_{\theta}(\kappa|x) - y(\kappa)| \quad (21)$$

We cannot apply Lemma B.6 directly to the Wasserstein loss class, because the Wasserstein loss is only defined on probability distributions, so 0 is not a valid input. To get around this problem, we assume the hypothesis space  $\mathcal{H}$  used in learning is of the form

$$\mathcal{H} = \{\mathfrak{s} \circ h^o : h^o \in \mathcal{H}^o\} \quad (22)$$

where  $\mathcal{H}^o$  is a function class that maps into  $\mathbb{R}^K$ , and  $\mathfrak{s}$  is the softmax function defined as  $\mathfrak{s}(o) = (\mathfrak{s}_1(o), \dots, \mathfrak{s}_K(o))$ , with

$$\mathfrak{s}_k(o) = \frac{e^{o_k}}{\sum_j e^{o_j}}, \quad k = 1, \dots, K \quad (23)$$

The softmax layer produce a valid probability distribution from arbitrary input, and this is consistent with commonly used models such as Logistic Regression and Neural Networks. By working with the log of the groundtruth labels, we can also add a softmax layer to the labels.

**Lemma B.9** (Proposition 2 of [28]). *The Wasserstein distances  $W_p(\cdot, \cdot)$  are metrics on the space of probability distributions of  $\mathcal{K}$ , for all  $1 \leq p \leq \infty$ .*

**Proposition B.10.** *The map  $\iota : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}$  defined by  $\iota(y, y') = W_1^1(\mathfrak{s}(y), \mathfrak{s}(y'))$  satisfies*

$$|\iota(y, y') - \iota(\bar{y}, \bar{y}')| \leq 4C_M \|(y, y') - (\bar{y}, \bar{y}')\|_2 \quad (24)$$

for any  $(y, y'), (\bar{y}, \bar{y}') \in \mathbb{R}^K \times \mathbb{R}^K$ . And  $\iota(0, 0) = 0$ .

*Proof.* For any  $(y, y'), (\bar{y}, \bar{y}') \in \mathbb{R}^K \times \mathbb{R}^K$ , by Lemma B.9, we can use triangle inequality on the Wasserstein loss,

$$|\iota(y, y') - \iota(\bar{y}, \bar{y}')| = |\iota(y, y') - \iota(\bar{y}, y') + \iota(\bar{y}, y') - \iota(\bar{y}, \bar{y}')| \leq \iota(y, \bar{y}) + \iota(y', \bar{y}')$$

Following Corollary B.8, it continues as

$$|\iota(y, y') - \iota(\bar{y}, \bar{y}')| \leq C_M (\|\mathfrak{s}(y) - \mathfrak{s}(\bar{y})\|_1 + \|\mathfrak{s}(y') - \mathfrak{s}(\bar{y}')\|_1) \quad (25)$$

Note for each  $k = 1, \dots, K$ , the gradient  $\nabla_y \mathfrak{s}_k$  satisfies

$$\|\nabla_y \mathfrak{s}_k\|_2 = \left\| \left( \frac{\partial \mathfrak{s}_k}{\partial y_j} \right)_{j=1}^K \right\|_2 = \left\| (\delta_{kj} \mathfrak{s}_k - \mathfrak{s}_k \mathfrak{s}_j)_{j=1}^K \right\|_2 = \sqrt{\mathfrak{s}_k^2 \sum_{j=1}^K \mathfrak{s}_j^2 + \mathfrak{s}_k^2 (1 - 2\mathfrak{s}_k)} \quad (26)$$

By mean value theorem,  $\exists \alpha \in [0, 1]$ , such that for  $y_\theta = \alpha y + (1 - \alpha)\bar{y}$ , it holds that

$$\|\mathfrak{s}(y) - \mathfrak{s}(\bar{y})\|_1 = \sum_{k=1}^K \left| \langle \nabla_y \mathfrak{s}_k |_{y=y_{\alpha k}}, y - \bar{y} \rangle \right| \leq \sum_{k=1}^K \|\nabla_y \mathfrak{s}_k |_{y=y_{\alpha k}}\|_2 \|y - \bar{y}\|_2 \leq 2\|y - \bar{y}\|_2$$

because by (26), and the fact that  $\sqrt{\sum_j \mathfrak{s}_j^2} \leq \sum_j \mathfrak{s}_j = 1$  and  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for  $a, b \geq 0$ , it holds

$$\begin{aligned} \sum_{k=1}^K \|\nabla_y \mathfrak{s}_k\|_2 &= \sum_{k: \mathfrak{s}_k \leq 1/2} \|\nabla_y \mathfrak{s}_k\|_2 + \sum_{k: \mathfrak{s}_k > 1/2} \|\nabla_y \mathfrak{s}_k\|_2 \\ &\leq \sum_{k: \mathfrak{s}_k \leq 1/2} (\mathfrak{s}_k + \mathfrak{s}_k \sqrt{1 - 2\mathfrak{s}_k}) + \sum_{k: \mathfrak{s}_k > 1/2} \mathfrak{s}_k \leq \sum_{k=1}^K 2\mathfrak{s}_k = 2 \end{aligned}$$

Similarly, we have  $\|\mathfrak{s}(y') - \mathfrak{s}(\bar{y}')\|_1 \leq 2\|y' - \bar{y}'\|_2$ , so from (25), we know

$$|\iota(y, y') - \iota(\bar{y}, \bar{y}')| \leq 2C_M (\|y - \bar{y}\|_2 + \|y' - \bar{y}'\|_2) \leq 2\sqrt{2}C_M (\|y - \bar{y}\|_2^2 + \|y' - \bar{y}'\|_2^2)^{1/2}$$

then (24) follows immediately. The second conclusion follows trivially as  $\mathfrak{s}$  maps the zero vector to a uniform distribution.  $\square$

*Proof of Theorem 5.1.* Consider the loss function space preceded with a softmax layer

$$\mathcal{L} = \{\iota_\theta : (x, y) \mapsto W_1^1(\mathfrak{s}(h_\theta^o(x)), \mathfrak{s}(y)) : h_\theta^o \in \mathcal{H}^o\}$$

We apply Lemma B.6 to the  $4C_M$ -Lipschitz continuous function  $\iota$  in Proposition B.10 and the function space

$$\underbrace{\mathcal{H}^o \times \dots \times \mathcal{H}^o}_{K \text{ copies}} \times \underbrace{\mathcal{I} \times \dots \times \mathcal{I}}_{K \text{ copies}}$$

with  $\mathcal{I}$  a singleton function space with only the identity map. It holds

$$\hat{\mathfrak{R}}_S(\mathcal{L}) \leq 8C_M \left( K\hat{\mathfrak{R}}_S(\mathcal{H}^o) + K\hat{\mathfrak{R}}_S(\mathcal{I}) \right) = 8KC_M\hat{\mathfrak{R}}_S(\mathcal{H}^o) \quad (27)$$

because for the identity map, and a sample  $S = (y_1, \dots, y_N)$ , we can calculate

$$\hat{\mathfrak{R}}_S(\mathcal{I}) = \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{I}} \frac{1}{N} \sum_{i=1}^N \sigma_i f(y_i) \right] = \mathbb{E}_\sigma \left[ \frac{1}{N} \sum_{i=1}^N \sigma_i y_i \right] = 0$$

The conclusion of the theorem follows by combining (27) with Theorem B.3 and Lemma B.1.  $\square$

## C Connection with other standard measures

### C.1 Connection with multiclass classification

*Proof of Proposition 5.2.* Given that the label is a “one-hot” vector  $y = \mathbb{e}_\kappa$ , the set of transport plans (3) degenerates. Specifically, the constraint  $T^\top \mathbf{1} = \mathbb{e}_\kappa$  means that only the  $\kappa$ -th column of  $T$  could be non-zero. Above that, the constraint  $T\mathbf{1} = h_{\hat{\theta}}(\cdot|x)$  ensures that the  $\kappa$ -th column of  $T$  actually equals to  $h_{\hat{\theta}}(\cdot|x)$ . In other words, the set  $\Pi(h_{\hat{\theta}}(\cdot|x), \mathbb{e}_\kappa)$  contains only one feasible transport plan, so (2) can be computed directly as

$$W_p^p(h_{\hat{\theta}}(\cdot|x), \mathbb{e}_\kappa) = \sum_{\kappa' \in \mathcal{K}} M_{\kappa', \kappa} h_{\hat{\theta}}(\kappa'|x) = \sum_{\kappa' \in \mathcal{K}} d_{\mathcal{K}}^p(\kappa', \kappa) h_{\hat{\theta}}(\kappa'|x)$$

Now let  $\hat{\kappa} = \operatorname{argmax}_{\kappa} h_{\hat{\theta}}(\kappa|x)$  be the prediction, we have

$$h_{\hat{\theta}}(\hat{\kappa}|x) = 1 - \sum_{\kappa \neq \hat{\kappa}} h_{\hat{\theta}}(\kappa|x) \geq 1 - \sum_{\kappa \neq \hat{\kappa}} h_{\hat{\theta}}(\hat{\kappa}|x) = 1 - (K-1)h_{\hat{\theta}}(\hat{\kappa}|x)$$

Therefore,  $h_{\hat{\theta}}(\hat{\kappa}|x) \geq 1/K$ , so

$$W_p^p(h_{\hat{\theta}}(\cdot|x), \mathbb{e}_\kappa) \geq d_{\mathcal{K}}^p(\hat{\kappa}, \kappa) h_{\hat{\theta}}(\hat{\kappa}|x) \geq d_{\mathcal{K}}^p(\hat{\kappa}, \kappa)/K$$

The conclusion follows by applying Theorem 5.1 with  $p = 1$ .  $\square$

### C.2 Connection with the total variation distance

The total variation distance between two distributions  $\mu$  and  $\nu$  is defined as  $\operatorname{TV}(\mu, \nu) = \sup_{A \subset \mathcal{K}} |\mu(A) - \nu(A)|$ . It can be shown that

$$\operatorname{TV}(\mu, \nu) = \frac{1}{2} \sum_{\kappa \in \mathcal{K}} |\mu(\kappa) - \nu(\kappa)| = 1 - \sum_{\kappa \in \mathcal{K}} \min(\mu(\kappa), \nu(\kappa)) \quad (28)$$

*Proof of Proposition 5.3.* In the case of 0-1 ground metric, the transport cost becomes

$$\langle T, M \rangle = \sum_{\kappa, \kappa'} T_{\kappa, \kappa'} M_{\kappa, \kappa'} = 1 - \sum_{\kappa} T_{\kappa, \kappa} \quad (29)$$

Moreover, by the constraint  $T \in \Pi(\mu, \nu)$ , we have

$$\mu(\kappa) = \sum_{\kappa'} T_{\kappa, \kappa'} = T_{\kappa, \kappa} + \sum_{\kappa' \neq \kappa} T_{\kappa, \kappa'} \geq T_{\kappa, \kappa}, \quad \nu(\kappa) = \sum_{\kappa'} T_{\kappa', \kappa} = T_{\kappa, \kappa} + \sum_{\kappa' \neq \kappa} T_{\kappa', \kappa} \geq T_{\kappa, \kappa}$$

Therefore, the minimum of (29) is achieved by  $T_{\kappa, \kappa}^* = \min(\mu(\kappa), \nu(\kappa))$  for the diagonal, with off-diagonal entries assigned arbitrarily as long as the constraints for  $\Pi(\mu, \nu)$  are met. As a result, it holds

$$W_{0-1}^1(\mu, \nu) = \langle T^*, M^{0-1} \rangle = 1 - \sum_{\kappa} \min(\mu(\kappa), \nu(\kappa)) \quad (30)$$

Following (28), we can see  $W_{0-1}^1$  equals to the total variation distance, which is also a scaled version of the  $\ell_1$  distance.  $\square$

### C.3 Connection with the Jaccard distance

Let  $W_{0-1}^1(\cdot, \cdot)$  be the Wasserstein distance under the 0-1 metric defined by  $d(\kappa, \kappa') = \mathbb{1}_{\kappa \neq \kappa'}$ , then we have the following characterization of the Wasserstein distance between two uniform distributions over regions.

**Lemma C.1.** *Let  $A, B \subset \mathcal{K}$ , then*

$$1 - W_{0-1}^1(\mathbb{U}^A, \mathbb{U}^B) = \frac{|A \cap B|}{\max(|A|, |B|)} \quad (31)$$

*Proof.* The overlapping mass that does not need to transport is given by

$$m_0 = \min\left(\frac{1}{|A|}, \frac{1}{|B|}\right) |A \cap B| \quad (32)$$

Since under 0-1 metric, any transport of mass has unit cost, so the minimum attainable transport cost is

$$1 \times (1 - m_0) = 1 - \frac{|A \cap B|}{\max(|A|, |B|)} \quad (33)$$

□

*Proof of Proposition 5.4.* Given that  $d_J(A, B) \leq \varepsilon$ , by Lemma C.1, it holds

$$W_{0-1}^1(\mathcal{U}^A, \mathcal{U}^B) = 1 - \frac{|A \cap B|}{\max(|A|, |B|)} \leq 1 - \frac{|A \cap B|}{|A \cup B|} = d_J(A, B) \leq \varepsilon$$

Conversely, given that  $W_{0-1}^1(\mathcal{U}^A, \mathcal{U}^B) \leq \varepsilon$ , again by Lemma C.1, it holds

$$\frac{|A \cap B|}{\max(|A|, |B|)} \geq 1 - \varepsilon \quad (34)$$

By symmetry, without loss of generality, assume  $|A| \geq |B|$ , then

$$|A \cap B| \geq (1 - \varepsilon)|A| \quad (35)$$

As a result, we have

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|} \geq \frac{(1 - \varepsilon)|A|}{|A| + |A| - (1 - \varepsilon)|A|} = \frac{1 - \varepsilon}{1 + \varepsilon}$$

The conclusion follows as

$$d_J(A, B) = 1 - J(A, B) \leq 1 - \frac{1 - \varepsilon}{1 + \varepsilon} = \frac{2\varepsilon}{1 + \varepsilon} \leq 2\varepsilon$$

□

*Proof of Corollary 5.5.* Note when the alphabet consists of integer coordinates of pixels, for any  $\kappa \neq \kappa'$ ,  $\|\kappa - \kappa'\|_p^p \geq 1$  for any  $p \geq 1$ . Therefore,  $M_{\kappa, \kappa'}^{0-1} \leq M_{\kappa, \kappa'}^p$ , i.e. the 0-1 ground metric matrix is bounded by the  $p$ -Euclidean ground metric matrix, elementwise. Let  $T^*$  be the optimal transport plan under the  $p$ -Euclidean ground metric, which is also a feasible transport plan under the 0-1 ground metric. So

$$W_{0-1}^1(\mathcal{U}^A, \mathcal{U}^B) = \inf_{T \in \Pi(\mathcal{U}^A, \mathcal{U}^B)} \langle T, M^{0-1} \rangle \leq \langle T^*, M^{0-1} \rangle \leq \langle T^*, M^p \rangle = W_p^p(\mathcal{U}^A, \mathcal{U}^B)$$

The conclusion follows by a direct application of Proposition 5.4. □

#### C.4 Relation to the Jaccard distance (non-uniform case)

**Lemma C.2.** Let  $A, B \subset \mathcal{K}$ , and  $\mu^A, \nu^B$  are two probability distributions supported on  $A, B$ , respectively, then

$$W_{0-1}^1(\mu^A, \nu^B) = 1 - \min\{\mu^A(A \cap B), \nu^B(A \cap B)\} \quad (36)$$

*Proof.* The amount of mass that completely matches is  $\min\{\mu^A(A \cap B), \nu^B(A \cap B)\}$ . The rest of the mass needs to be moved around with unit cost 1, so the total minimum transport cost is  $1 - \min\{\mu^A(A \cap B), \nu^B(A \cap B)\}$ . □



**Proposition C.3.** Let  $\mu^A$  and  $\nu^B$  be two probability measures supported on  $A$  and  $B$ , respectively. Denote

$$\mu^* = \max_{\kappa \in A} \mu^A(\kappa), \quad \mu_o = \min_{\kappa \in A} \mu^A(\kappa)$$

$$\nu^* = \max_{\kappa \in B} \nu^B(\kappa), \quad \nu_o = \min_{\kappa \in B} \nu^B(\kappa)$$

then  $W_{0-1}^1(\mu^A, \nu^B) \leq \varepsilon$  implies

$$d_J(A, B) \leq \frac{2C^* - 2C_o(1 - \varepsilon)}{2C^* - (1 - \varepsilon)C_o} \quad (37)$$

where  $C^* \geq \max(\mu^*, \nu^*)$  and  $0 < C_o \leq \min(\mu_o, \nu_o)$ .

*Proof.* Notice obviously, for any  $X \subset \mathcal{K}$ , we have the following properties

$$\mu_o|X| \leq \mu^A(X) \leq \mu^*|X| \quad (38)$$

$$\nu_o|X| \leq \nu^B(X) \leq \nu^*|X| \quad (39)$$

$$(40)$$

Let us first assume that  $W_{0-1}^1(\mu^A, \nu^B) \leq \varepsilon$ , following Lemma C.2, it implies

$$1 - W_{0-1}^1(\mu^A, \nu^B) = \min(\mu^A(A \cap B), \nu^B(A \cap B)) \geq 1 - \varepsilon \quad (41)$$

On the other hand,

$$d_J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

So we can get an upper bound of  $d_J(A, B)$  by deriving a lower bound on  $|A \cap B|$  and an upper bound on  $|A \cup B|$ . By (38), (39), and (41), it holds

$$|A \cap B| \geq \max \left\{ \frac{\mu^A(A \cap B)}{\mu^*}, \frac{\nu^B(A \cap B)}{\nu^*} \right\} \geq \max \left\{ \frac{1 - \varepsilon}{\mu^*}, \frac{1 - \varepsilon}{\nu^*} \right\} \geq \frac{1 - \varepsilon}{C^*}$$

where  $C^* \geq \max\{\mu^*, \nu^*\}$ . Similarly, we have

$$|A \cup B| = |A| + |B| - |A \cap B| \leq \frac{1}{\mu_o} + \frac{1}{\nu_o} - \frac{1 - \varepsilon}{C^*} \leq \frac{2}{C_o} - \frac{1 - \varepsilon}{C^*}$$

where  $0 < C_o \leq \min\{\mu_o, \nu_o\}$ . It then follows that

$$d_J(A, B) \leq 1 - \frac{\frac{1 - \varepsilon}{C^*}}{\frac{2}{C_o} - \frac{1 - \varepsilon}{C^*}} \leq 1 - \frac{C_o(1 - \varepsilon)}{2C^* - (1 - \varepsilon)C_o} \leq \frac{2C^* - 2C_o(1 - \varepsilon)}{2C^* - (1 - \varepsilon)C_o}$$

□

**Proposition C.4.** Following the same notation of Proposition C.3,  $d_J(A, B) \leq \varepsilon$  implies

$$W_{0-1}^1(\mu^A, \nu^B) \leq \frac{2(C^* - C_o) + \varepsilon(2C_o - C^*)}{C^*(2 - \varepsilon)} \quad (42)$$

*Proof.* By Lemma C.2, in order to upper bound  $W_{0-1}^1(\mu^A, \nu^B)$ , we only need to derive lower bounds for  $\mu^A(A \cap B)$  and  $\nu^B(A \cap B)$ . By  $d_J(A, B) \leq \varepsilon$ , it holds

$$1 - \varepsilon \leq \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|} \leq \frac{|A \cap B|}{1/\mu^* + 1/\nu^* - |A \cap B|}$$

where the inequality is from (38) and (39). As a result,

$$|A \cap B| \geq \frac{1 - \varepsilon}{2 - \varepsilon} \left( \frac{1}{\mu^*} + \frac{1}{\nu^*} \right) \geq \frac{1 - \varepsilon}{2 - \varepsilon} \frac{2}{C^*}$$

where  $C^* \geq \max\{\mu^*, \nu^*\}$ . By (38) again, we get

$$\mu^A(A \cap B) \geq \mu_o|A \cap B| \geq \frac{1 - \varepsilon}{2 - \varepsilon} \frac{2\mu_o}{C^*}$$

Similarly, we have

$$\nu^B(A \cap B) \geq \frac{1 - \varepsilon}{2 - \varepsilon} \frac{2\nu_o}{C^*}$$

Combining, we get

$$1 - W_{0-1}^1(\mu^A, \nu^B) = \min \{ \mu^A(A \cap B), \nu^B(A \cap B) \} \geq \frac{2C_o}{C^*} \frac{1 - \varepsilon}{2 - \varepsilon}$$

where  $0 < C_o \leq \min\{\mu_o, \nu_o\}$ . The conclusion follows immediately.  $\square$

REMARK: For the case of uniform distributions,  $C^* = C_o$ , Proposition C.3 and Proposition C.4 fall back to Proposition 5.4.

## D Empirical study

### D.1 Noisy label example

We illustrate the phenomenon of semantic label noise of human labelers with a synthetic example. Consider a multiclass classification problem, where the labels corresponds to the vertices on a  $D \times D$  lattice on the 2D plane. The Euclidean distance in  $\mathbb{R}^2$  is used to measure the semantic similarity between labels. The observations for each category are samples of a isotropic Gaussian distribution centered at the corresponding vertex. Given a noise level  $t$ , we choose to flip the label of each training sample to one of the neighboring categories<sup>5</sup> with probability  $t$ . Figure 8 shows the training set for  $3 \times 3$  lattice with noise level  $t = 0.1$  and  $t = 0.5$ , respectively.

We repeat 10 times for noise levels  $t = 0.1, 0.2, \dots, 0.9$  and  $D = 3, 4, \dots, 7$ . A multiclass classifier based on logistic regression is trained with the standard divergenced based loss<sup>6</sup> and the proposed Wasserstein loss<sup>7</sup>, respectively. The performance is measured by the Euclidean distance between the predicted class and the true class, averaged on the test set. Figure 2 compares the performance of the two loss functions.

### D.2 Full figure for the MNIST example

The full version of Figure 4 from Section 6.1 is shown in Figure 9.

### D.3 Details of the Flickr tag prediction experiment

From the tags in the Yahoo Flickr Creative Commons dataset, we filtered out those not occurring in the WordNet<sup>8</sup> database, as well those whose dominant lexical category was "noun.location" or "noun.time." We also filtered out by hand nouns referring to geographical location or nationality, proper nouns, numbers, photography-specific vocabulary, and several words not generally descriptive of visual content (such as "annual" and "demo"). From the remainder, the 1000 most frequently occurring tags were used.

We list some of the 1000 selected tags here. The 50 most frequently occurring tags: *travel, square, wedding, art, flower, music, nature, party, beach, family, people, food, tree, summer, water, concert, winter, sky, snow, street, portrait, architecture, car, live, trip, friend, cat, sign, garden, mountain, bird, sport, light, museum, animal, rock, show, spring, dog, film, blue, green, road, girl, event, red, fun, building, new, cloud.* ... and the 50 least frequent tags: *arboretum, chick, sightseeing, vineyard, animalia, burlesque, key, flat, whale, swiss, giraffe, floor, peak, contemporary, scooter, society, actor, tomb, fabric, gala, coral, sleeping, lizard, performer, album, body, crew, bathroom, bed, cricket, piano, base, poetry, master, renovation, step, ghost, freight, champion, cartoon, jumping, crochet, gaming, shooting, animation, carving, rocket, infant, drift, hope.*

<sup>5</sup>Connected vertices on the lattice are considered neighbors, and the Euclidean distance between neighbors is set to 1.

<sup>6</sup>Corresponds to maximum likelihood estimation of the logistic regression model.

<sup>7</sup>In this special case, corresponds to weighted maximum likelihood estimation, c.f. Section C.1.

<sup>8</sup><http://wordnet.princeton.edu>

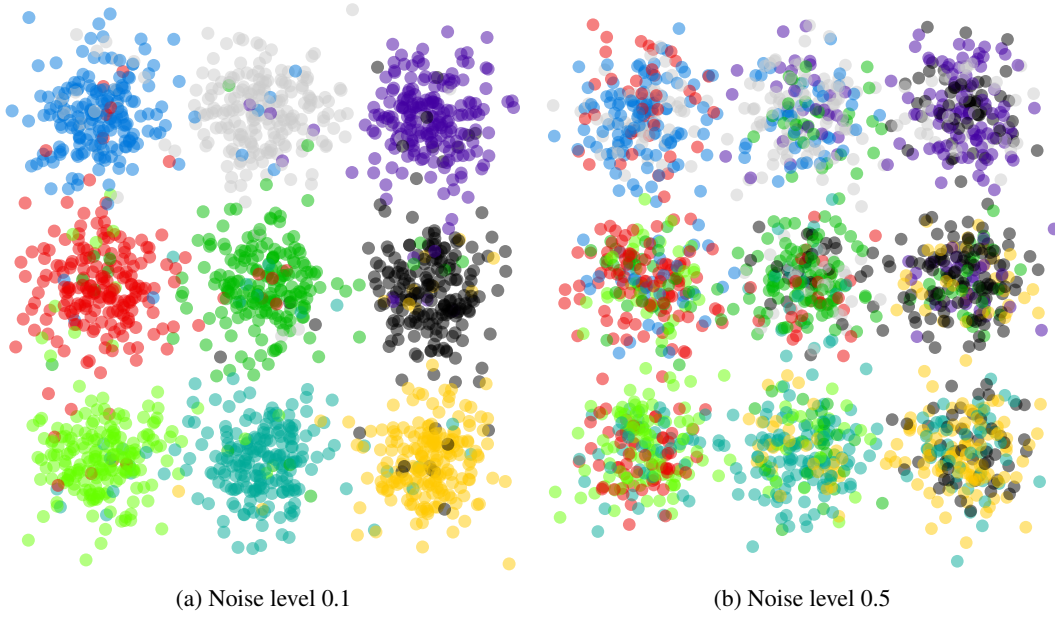


Figure 8: Illustration of training samples on a 3x3 lattice with different noise levels.

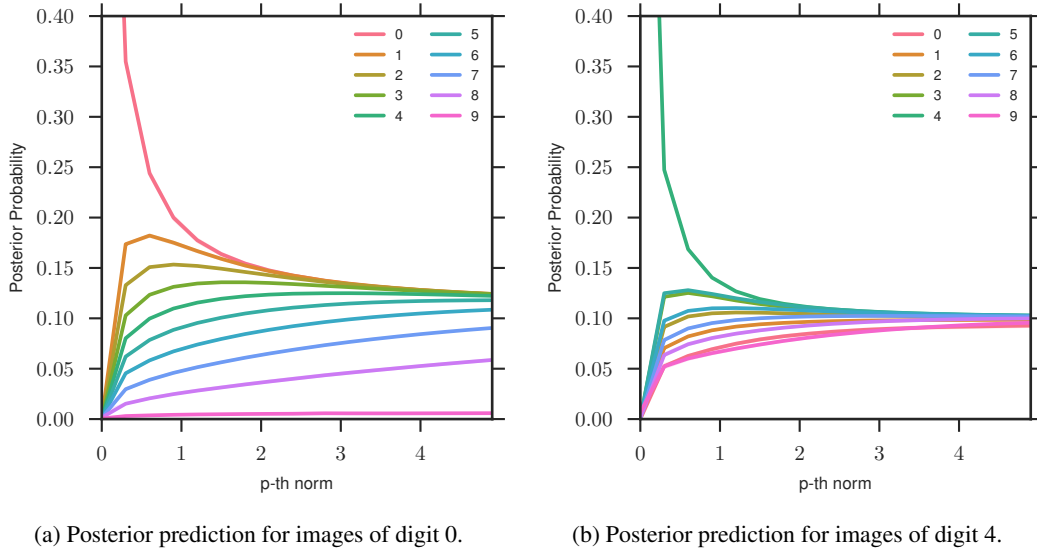


Figure 9: Each curve is the predicted probability for a target digit from models trained with different  $p$  values for the ground metric.

We train a multiclass linear logistic regression model with a linear combination of the Wasserstein loss and the KL divergence-based loss. The Wasserstein loss between the prediction and the normalized groundtruth is computed with an entropic regularization formulation using 10 iterations of the Sinkhorn-Knopp algorithm. Based on our observation of the ground metric matrix, we use  $p$ -norm with  $p = 13$ , and set  $\lambda = 50$ . This makes sure that the matrix  $\mathbf{K}$  is reasonable sparse, enforcing semantic smoothness only in each local neighborhood. Stochastic gradient descent with a mini-batch size of 100, and momentum 0.7 is run for 100,000 iterations to optimize the objective function on the training set. The baseline is trained under the same setting, but with only the KL loss function.

To create the dataset with reduced redundancy, for each image in the training set, we compute the pairwise semantic distance for the groundtruth tags, and cluster them into “equivalent” tag-sets with a threshold of semantic distance 1.3. Within each tag-set, one random tag is selected.

Figure 10 shows some more test images and predictions randomly picked from the test set.



(a) **Flickr user tags:** zoo, run, mark; **our proposals:** running, summer, fun; **baseline proposals:** running, country, lake.



(b) **Flickr user tags:** travel, architecture, tourism; **our proposals:** sky, roof, building; **baseline proposals:** art, sky, beach.



(c) **Flickr user tags:** spring, race, training; **our proposals:** road, bike, trail; **baseline proposals:** dog, surf, bike.



(d) **Flickr user tags:** family, trip, house; **our proposals:** family, girl, green; **baseline proposals:** woman, tree, family.



(e) **Flickr user tags:** education, weather, cow, agriculture; **our proposals:** girl, people, animal, play; **baseline proposals:** concert, statue, pretty, girl.



(f) **Flickr user tags:** garden, table, gardening; **our proposals:** garden, spring, plant; **baseline proposals:** garden, decoration, plant.



(g) **Flickr user tags:** nature, bird, rescue; **our proposals:** bird, nature, wildlife; **baseline proposals:** ature, bird, baby.

Figure 10: Examples of images in the Flickr dataset. We show the groundtruth tags and as well as tags proposed by our algorithm and baseline.